VIS 2016

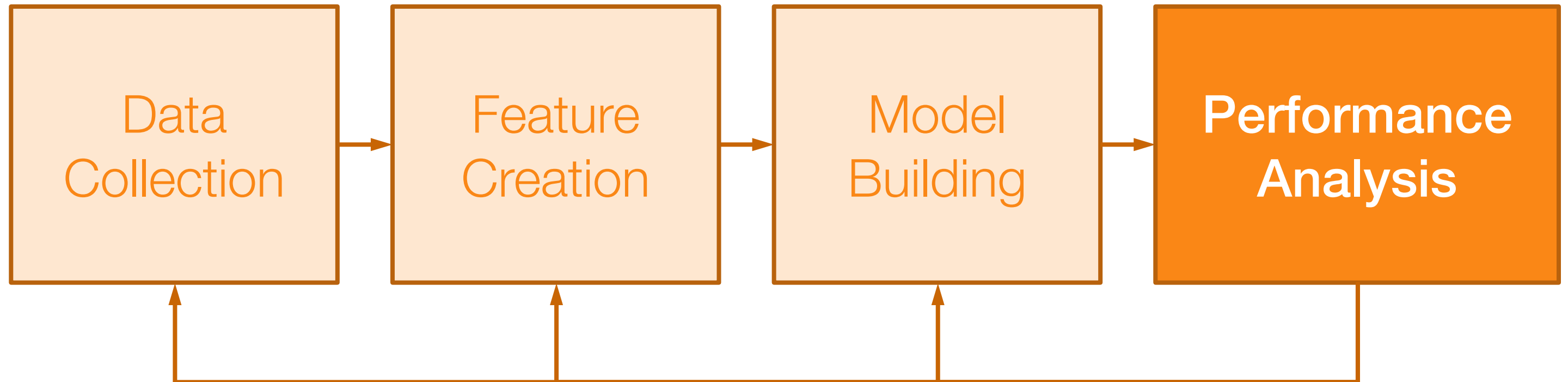# Squares: Supporting Interactive Performance Analysis for Multiclass Classifiers

**Donghao Ren**[1,2], Saleema Amershi[2], Bongshin Lee[2], Jina Suh[2] and Jason D. Williams[2]
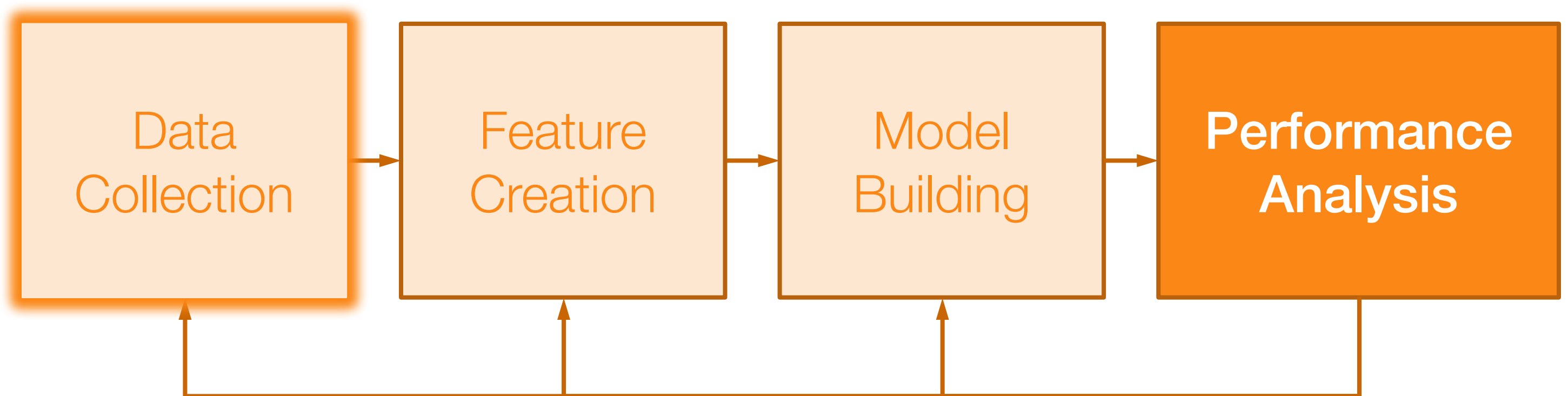
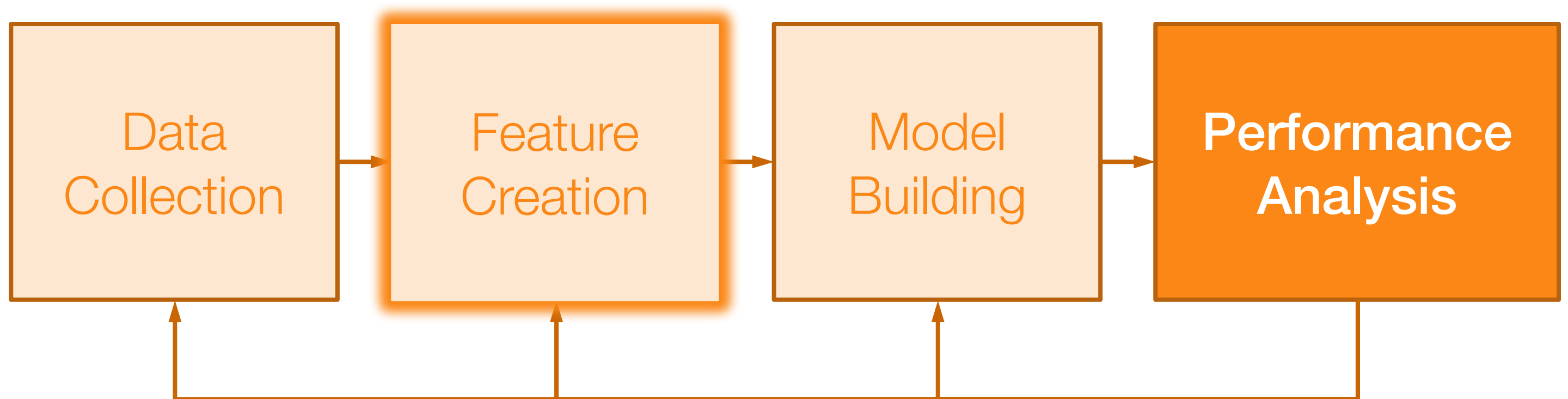[1] University of California, Santa Barbara
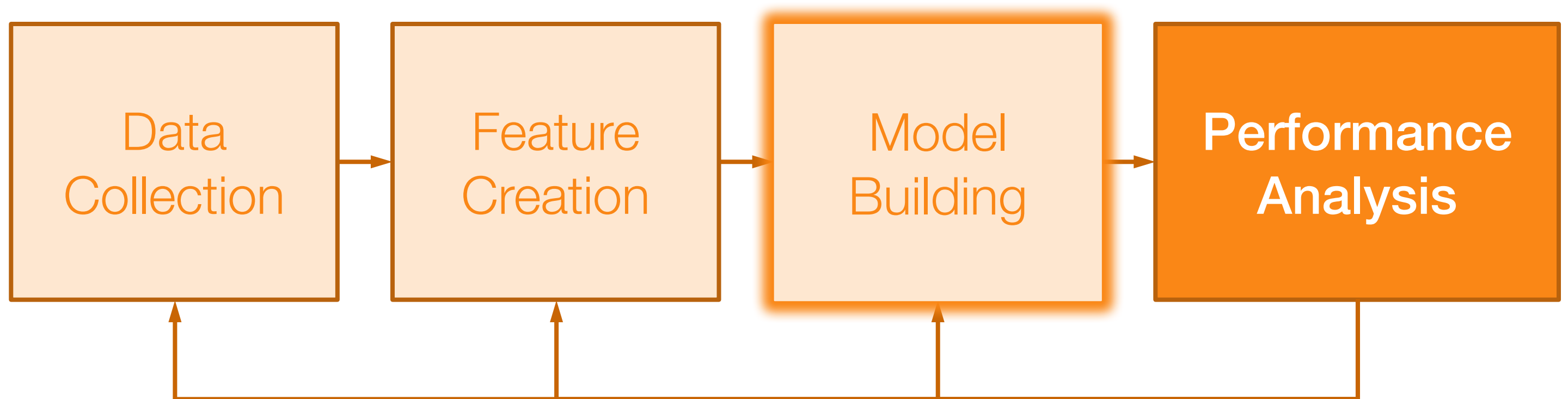[2] Microsoft Research, Redmond

# Performance analysis is critical in machine learning

# Performance analysis is critical in machine learning

# Performance analysis is critical in machine learning

# Performance analysis is critical in machine learning

# Common ways of performance analysis

- Summary statistics
  - Accuracy
  - Precision
  - Recall
  - Log-Loss
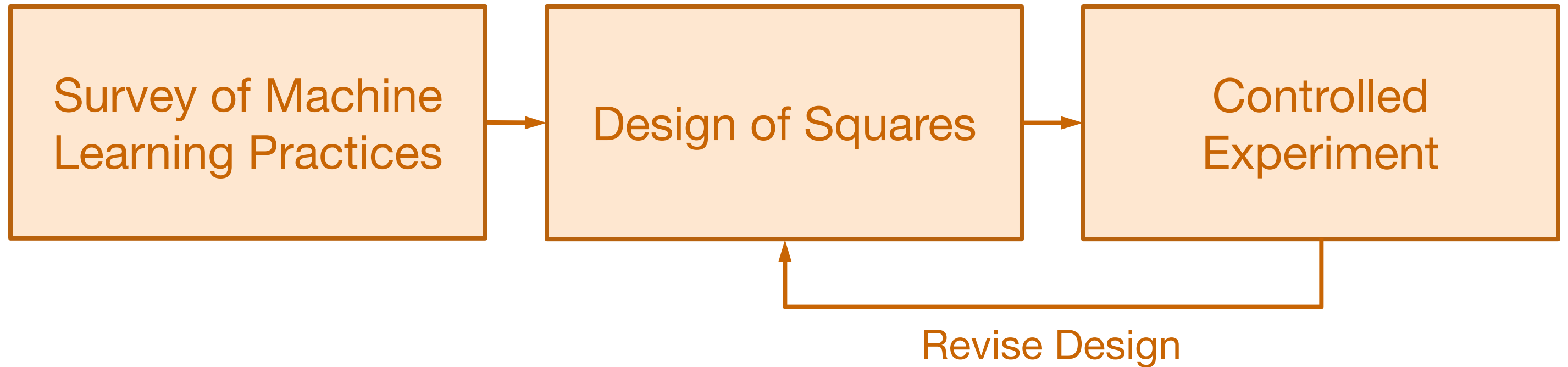  - …

- Confusion Matrix

Predicted Class



Actual Class

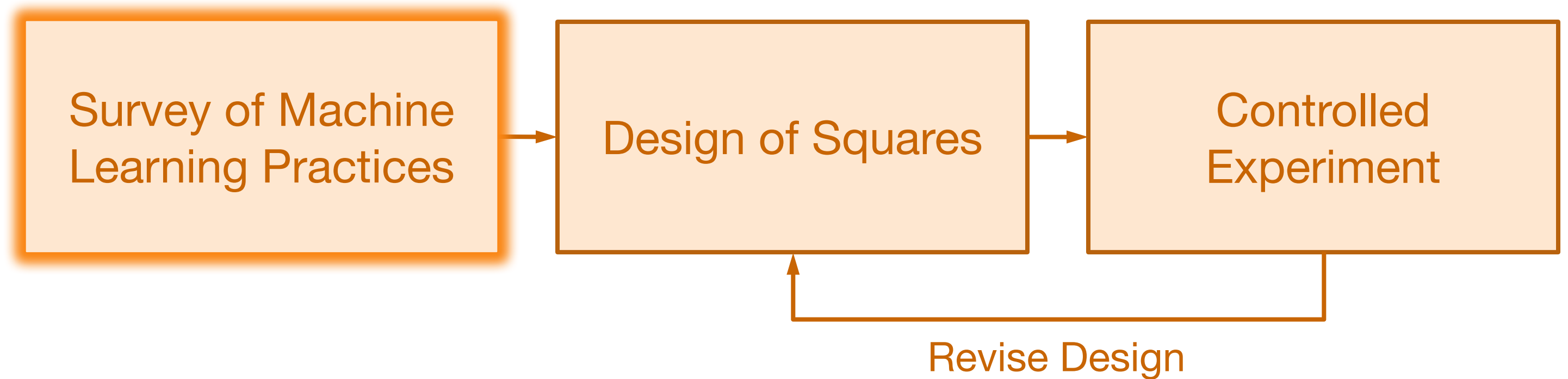| 89.3% | 0.6% | 6.9% | 3.3% |
| 31.4% | 42.0% | 19.4% | 7.2% |
| 18.6% | 0.4% | 79.8% | 1.2% |
| 16.3% | 1.1% | 2.4% | 80.2% |

# Problems

- Disconnected from the underlying data.

- Hide important information such as score distribution.

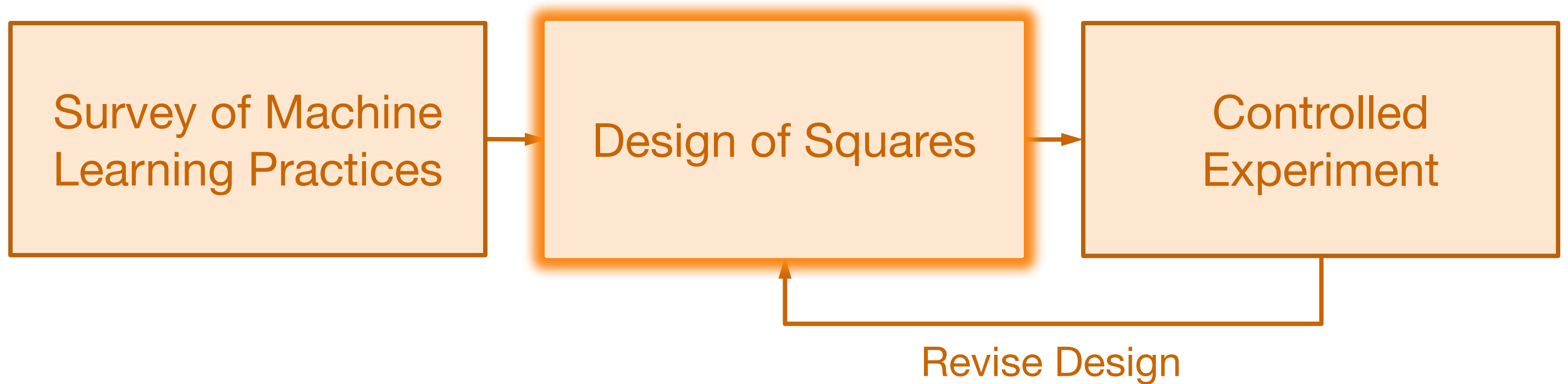- Not trivial to support **_multiclass_** classifiers.

# Squares

# Design Process


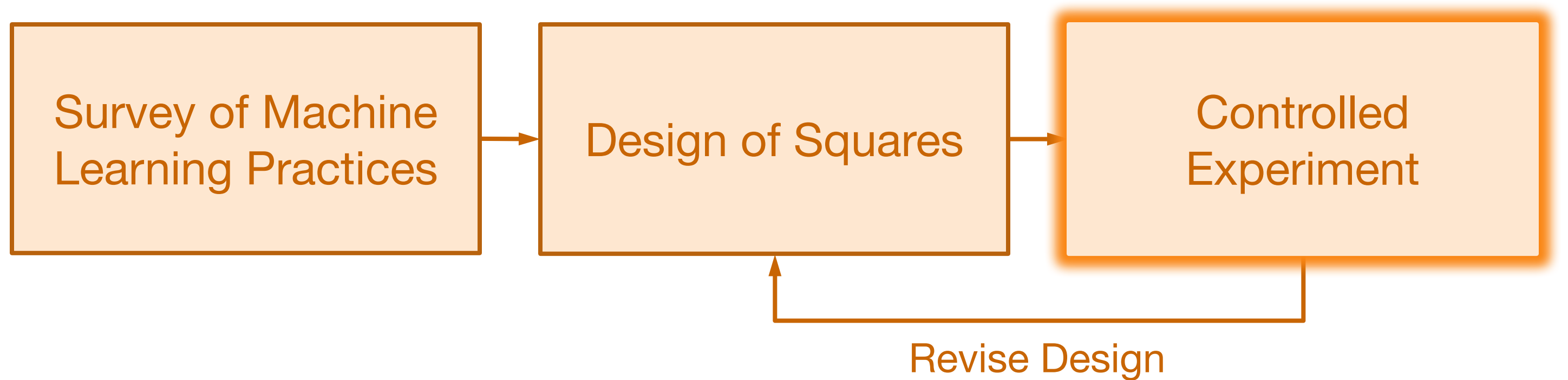
Survey of Machine Learning Practices → Design of Squares → Controlled Experiment

Revise Design

# Design Process

# Design Process

# Design Process

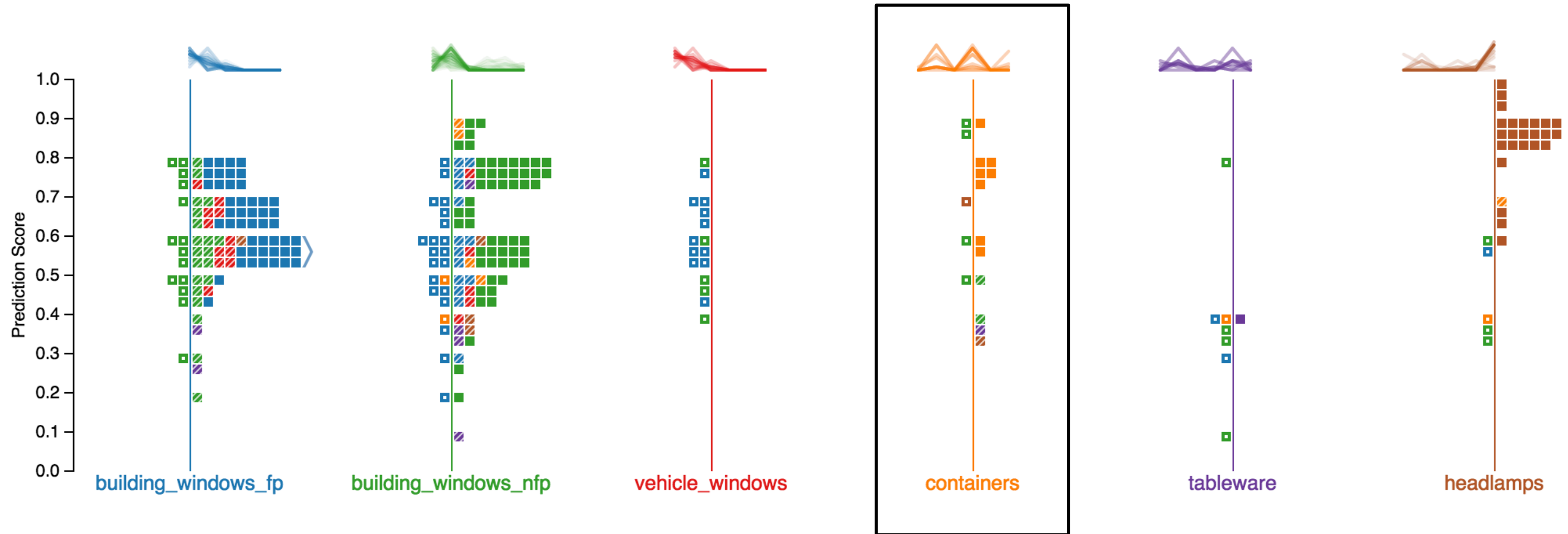Survey of Machine Learning Practices → Design of Squares → Controlled Experiment
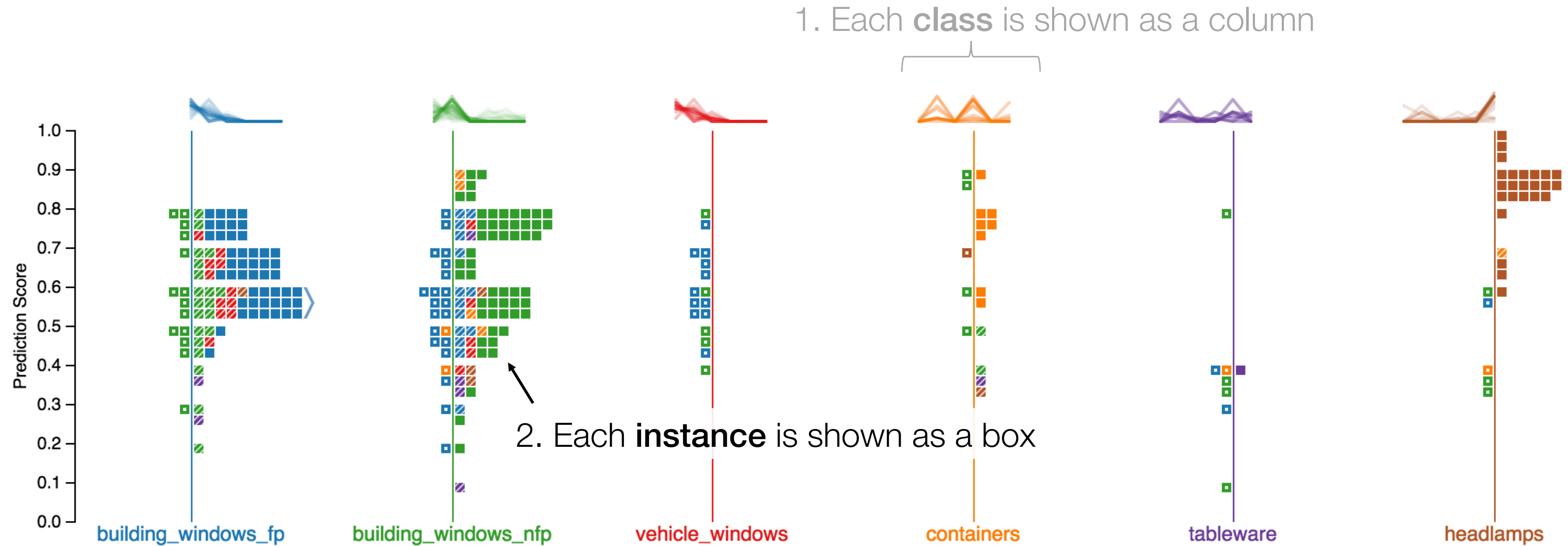
Revise Design

# Design Goals

- G1: Show performance at multiple levels of detail to help practitioners prioritize efforts.

  - Overall / Class-level / Instance-level

  - Error severity (errors with higher score on the wrong class are more severe)

- G2: Be agnostic to common performance metrics.

  - Support a wider range of scenarios.

- G3: Connect performance to data.

  - Provide access to data. Use small visual footprint to reserve space for scenario-dependent data access views.

# Squares Visualization Design

1. Each **class** is shown as a column

# Visualization Design



1. Each **class** is shown as a column

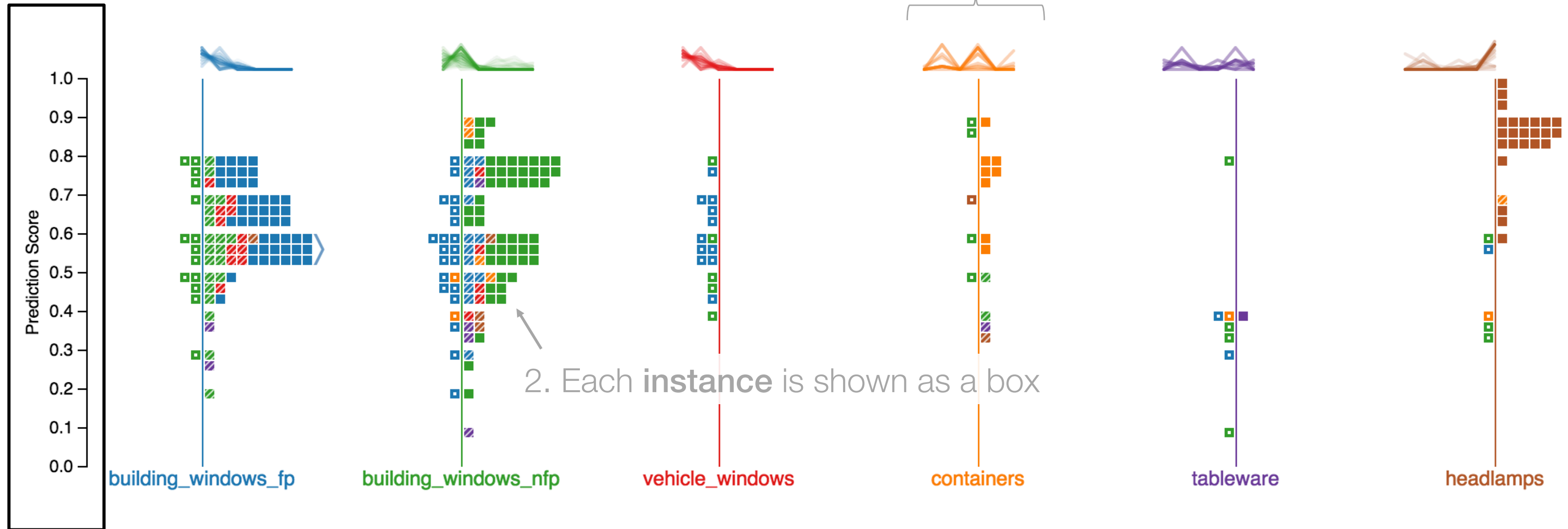2. Each **instance** is shown as a box
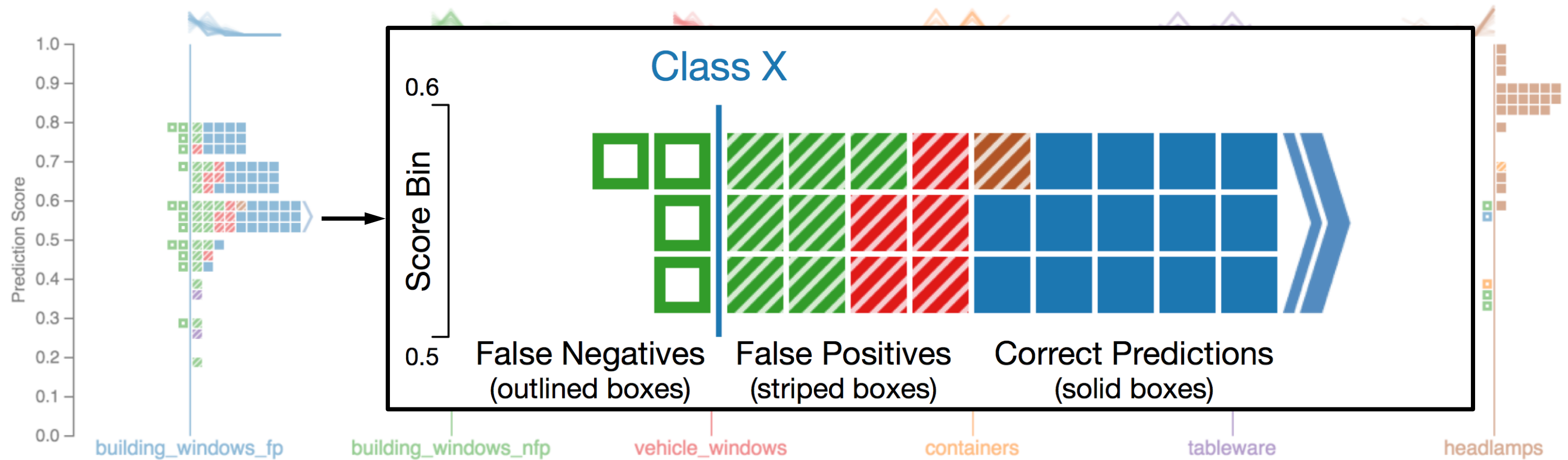
# Visualization Design



1. Each **class** is shown as a column

2. Each **instance** is shown as a box

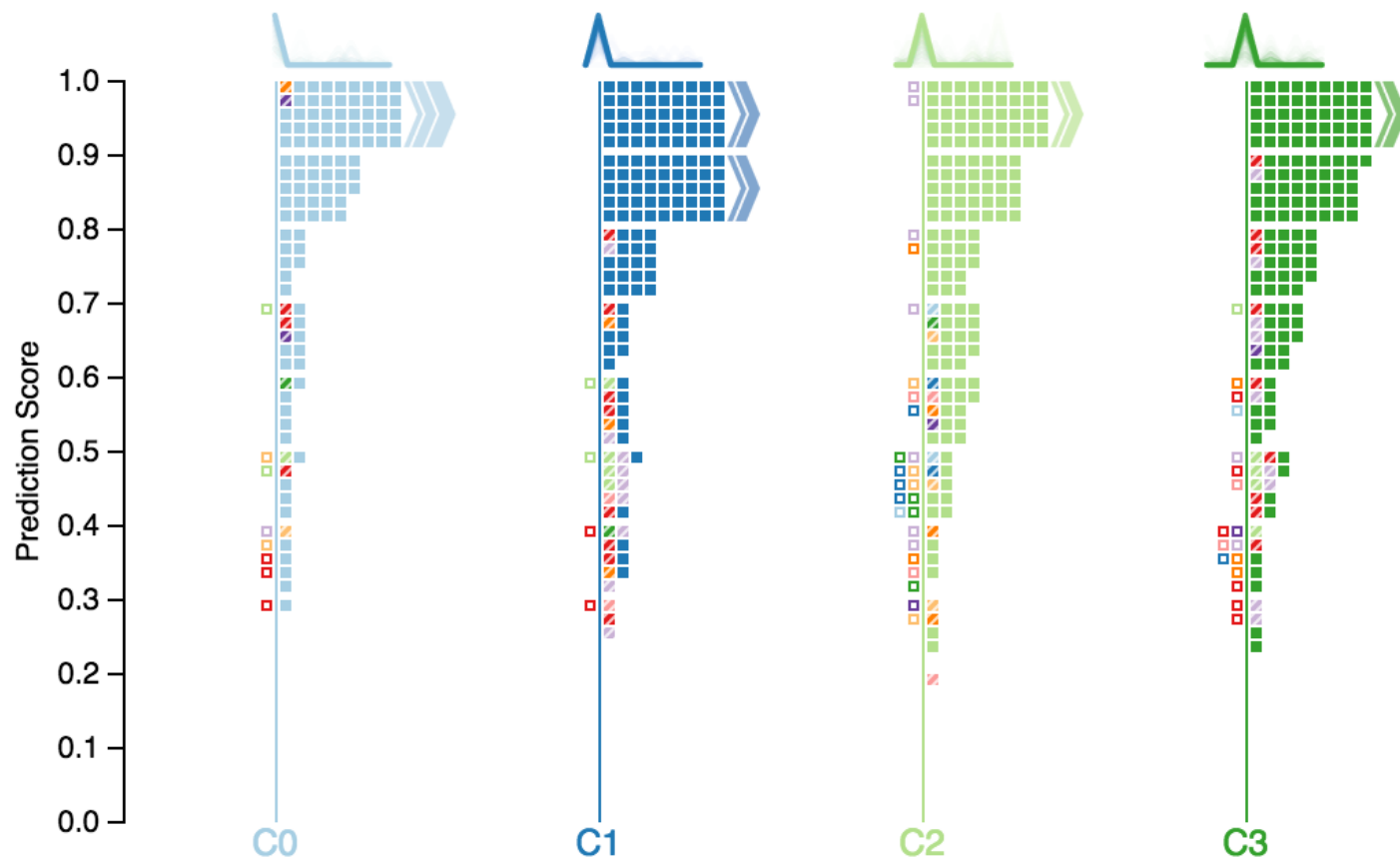3. Instances are binned according to **prediction scores**

# Visualization Design



Class X

Score Bin: 0.6 to 0.5

False Negatives (outlined boxes)
False Positives (striped boxes)
Correct Predictions (solid boxes)

Prediction Score axis: 0.0 to 1.0

building_windows_fp · building_windows_nfp · vehicle_windows · containers · tableware · headlamps
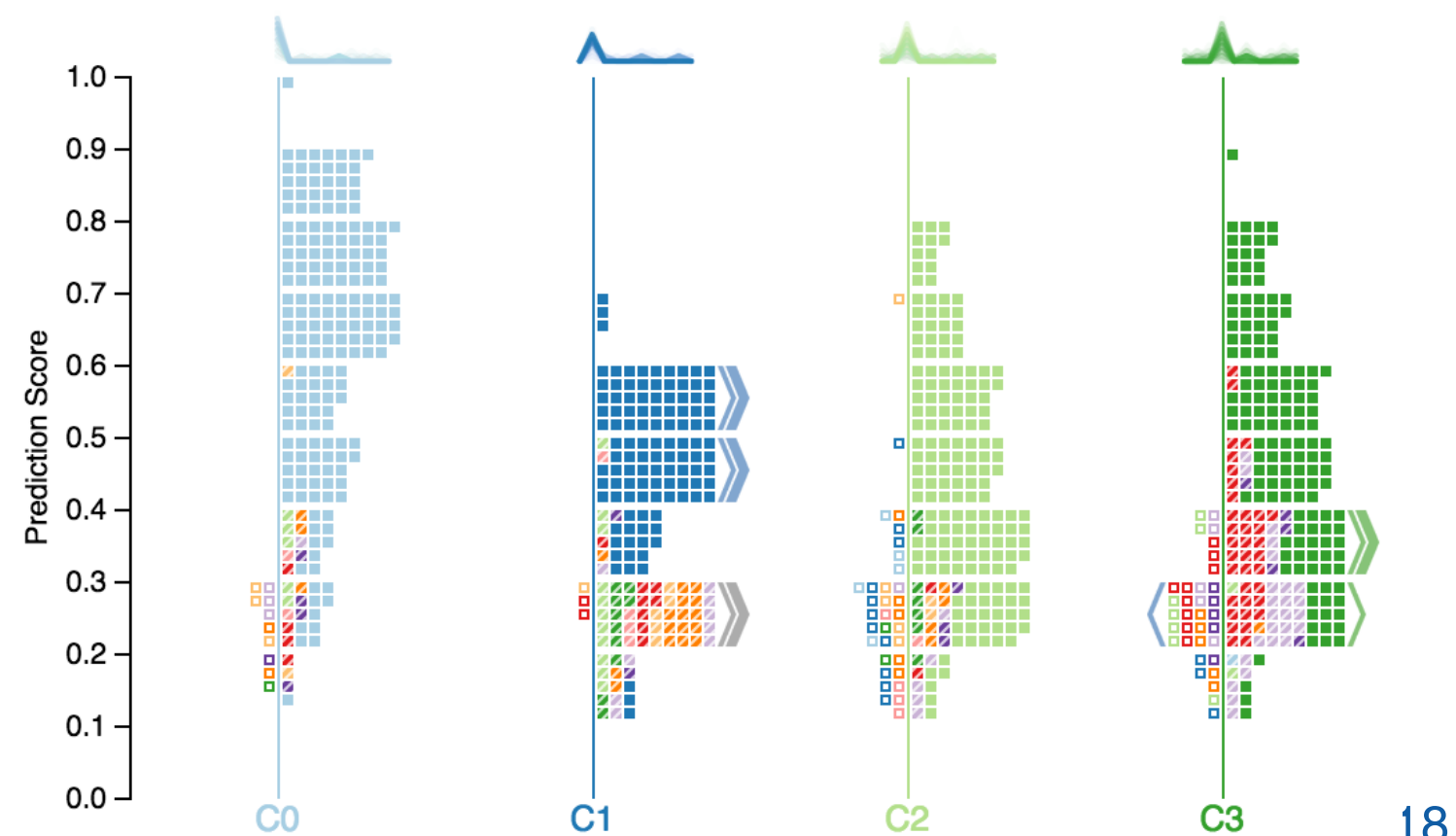
# Visualizing Count-Based Metrics: Overall Accuracy

- Accuracy: $\dfrac{\text{Correct Predictions}}{\text{Total \# of Instances}} = \dfrac{\blacksquare}{\blacksquare + \diagdown\diagup}$



Higher Accuracy

Lower Accuracy

# Visualizing Count-Based Metrics: Class-Level

- Class-level precision and recall:

Precision:

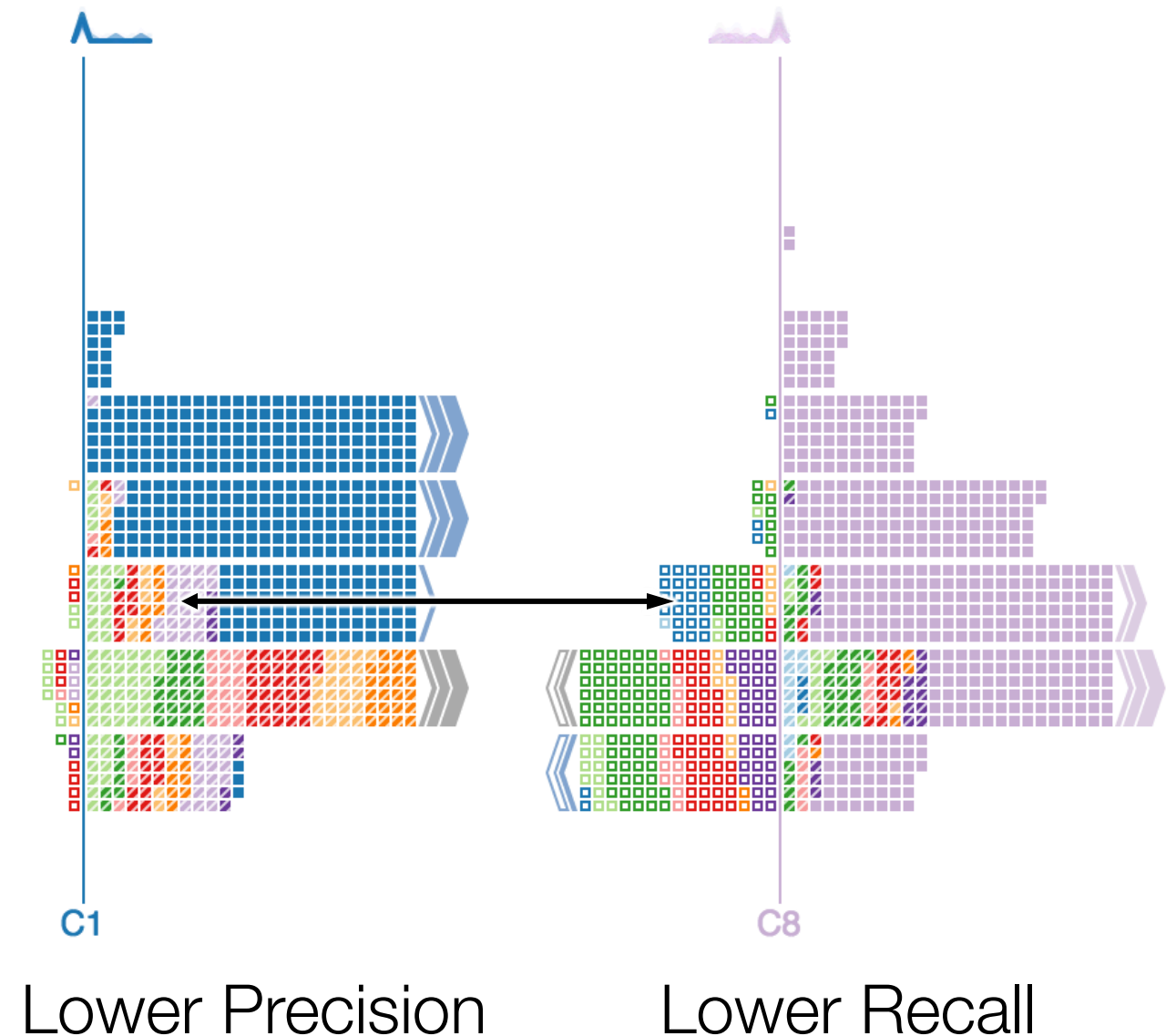$$\frac{TP}{TP + FP} = \frac{\blacksquare}{\blacksquare + \boxtimes}$$
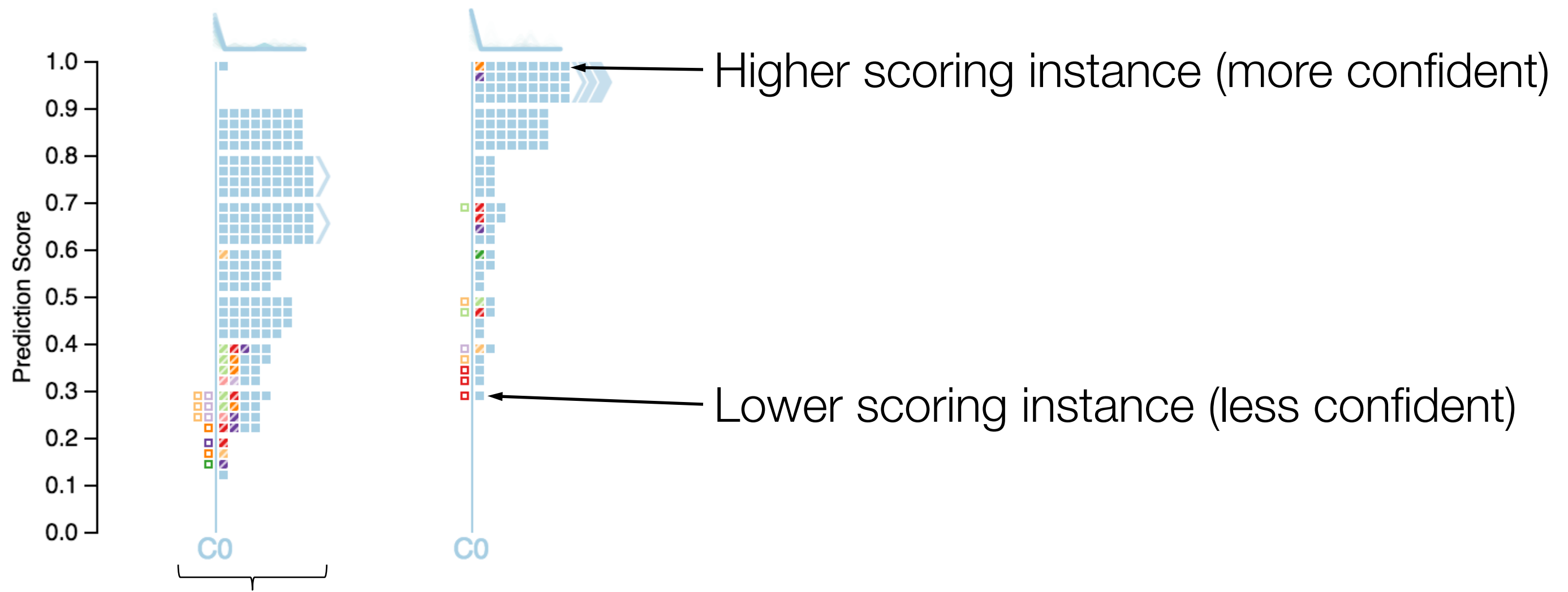
Recall:

$$\frac{TP}{TP + FN} = \frac{\blacksquare}{\blacksquare + \square}$$

FPs and FNs are comparably salient:

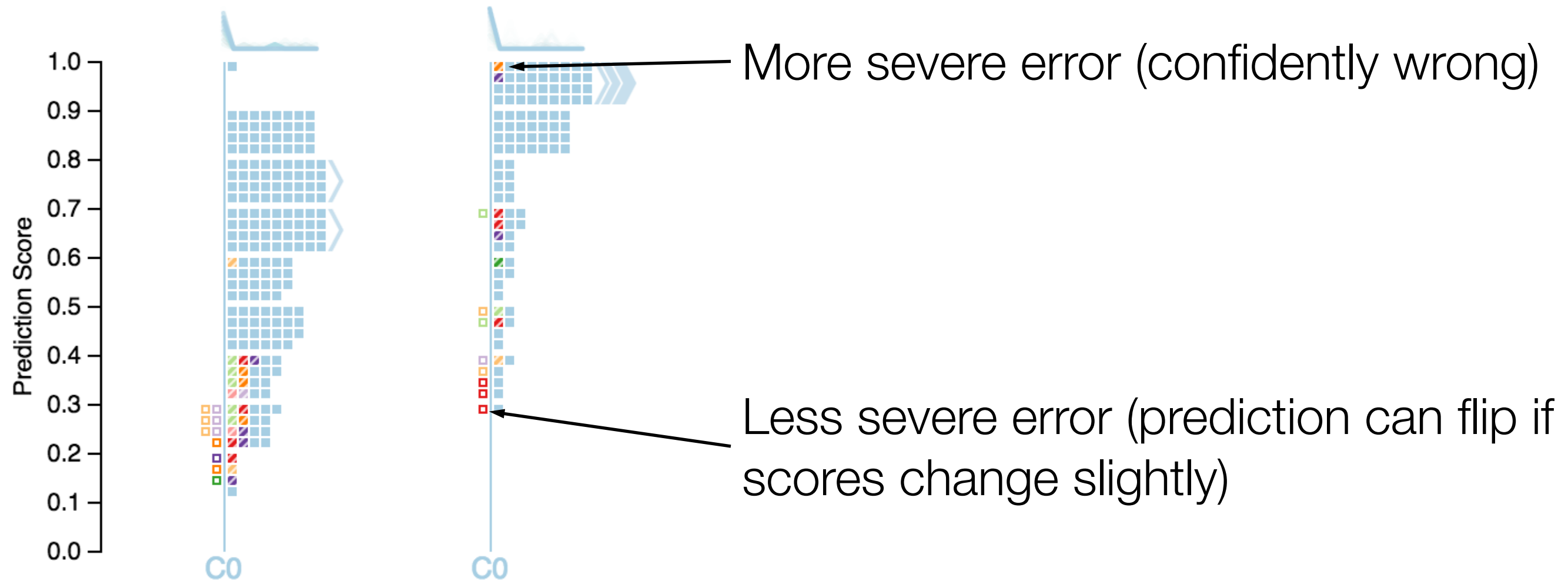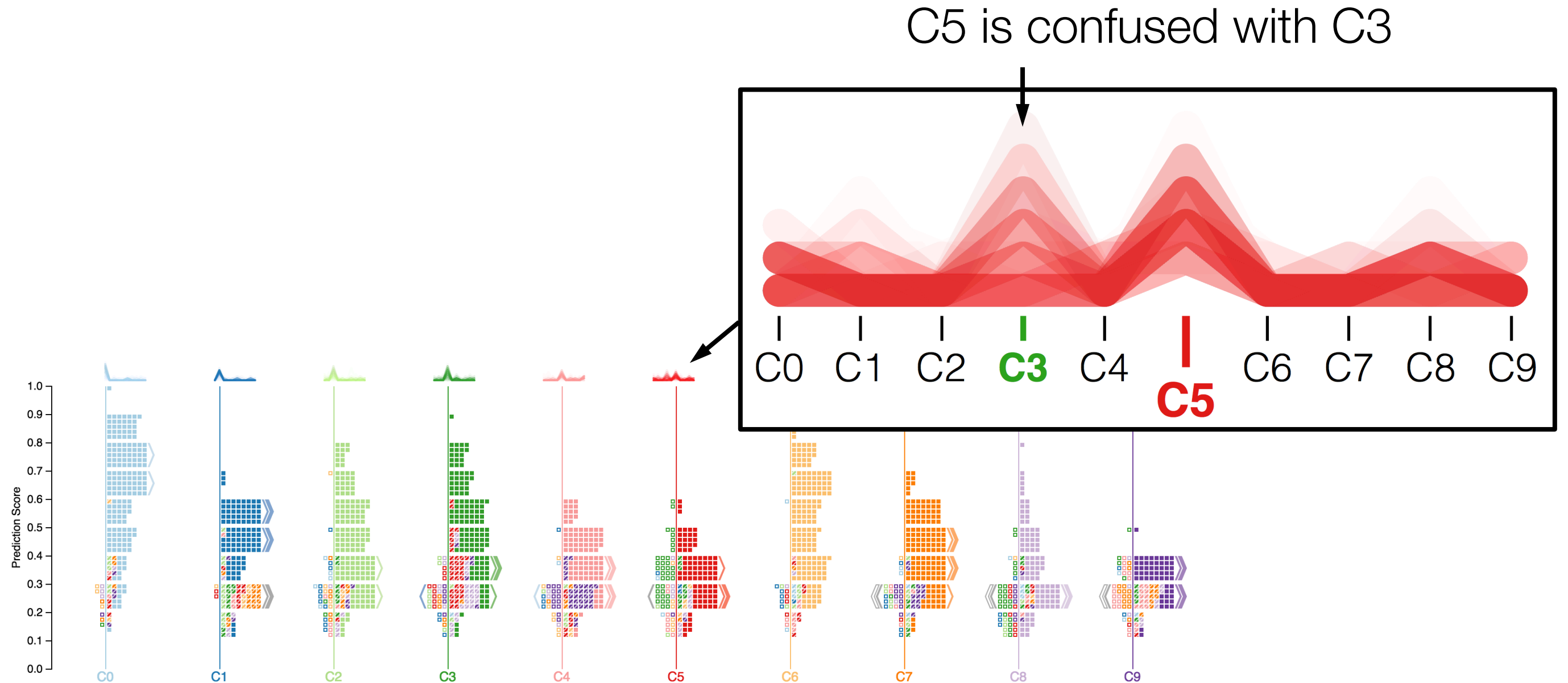One-to-one correspondence between outlined boxes and striped boxes

Lower Precision

Lower Recall

# Visualizing Score-Based Metrics



Higher scoring instance (more confident)

Lower scoring instance (less confident)

Worse score distribution

# Help Prioritizing Debugging Efforts



More severe error (confidently wrong)

Less severe error (prediction can flip if scores change slightly)

# Visualizing Confusion Between Classes



C5 is confused with C3

Dataset: **MNIST Handwritten Digits**   22

# Instance-Level Details



On-hover parallel coordinates for detailed scores

# Scalability



Each strip represents 10 boxes

Truncation indicators

Boxes

Strips

Stacks

# Scalability



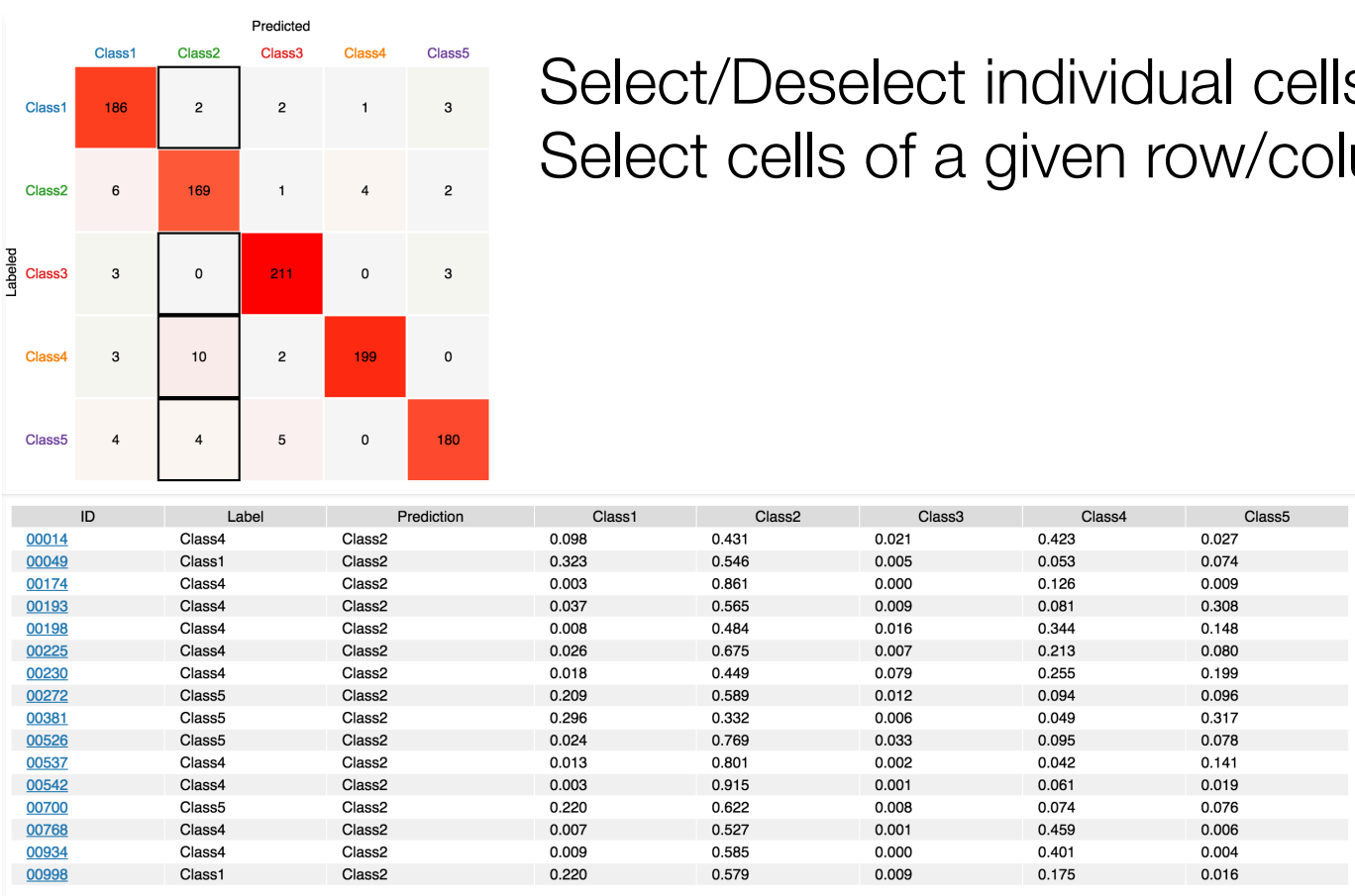Toggle between 3-levels of aggregation

Evaluation

# Controlled Experiment

- 24 participants

- Part 1: Comparison
  - Compare Squares against a commonly used ConfusionMatrix
  - Within-subject design

- Part 2: (Squares Only) Score Distribution
  - Evaluate Squares' ability to convey score distribution

# Part 1: Squares vs. Confusion Matrix



Select/Deselect individual cells.
Select cells of a given row/column.
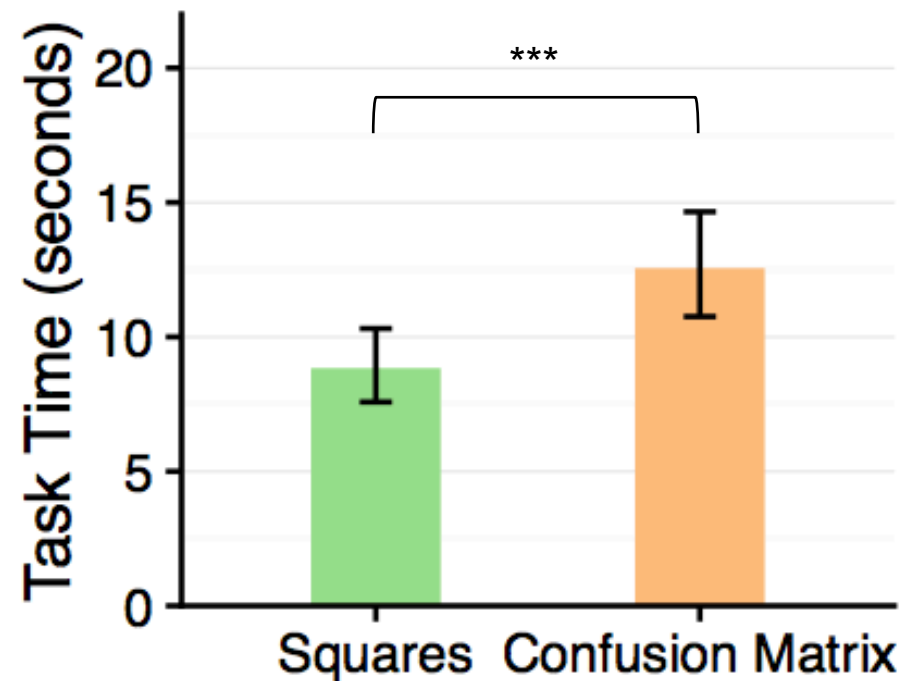
**Squares** with a Sortable Table

**Confusion Matrix** with a Sortable Table

# Part 1: Tasks

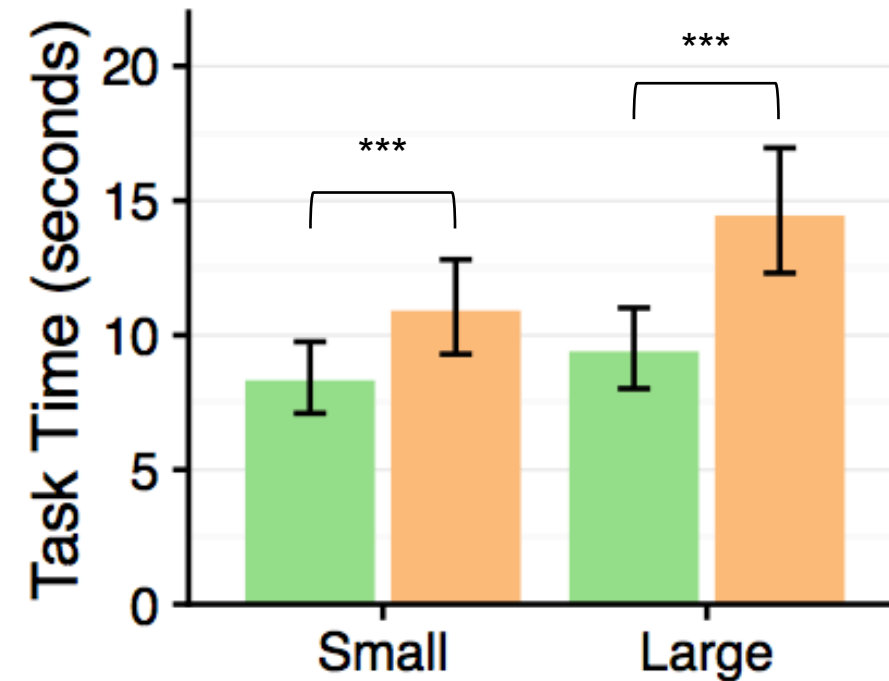- ## T1 – Overall
  - Select the classifier with the larger number of errors

- ## T2 – Class-level
  - Select one of the two classes with the most errors

- ## T3 – Instance-level
  - Select an error with a score of .9 or above in the wrong class

# Part 1: Squares Performed Better

- Task Time



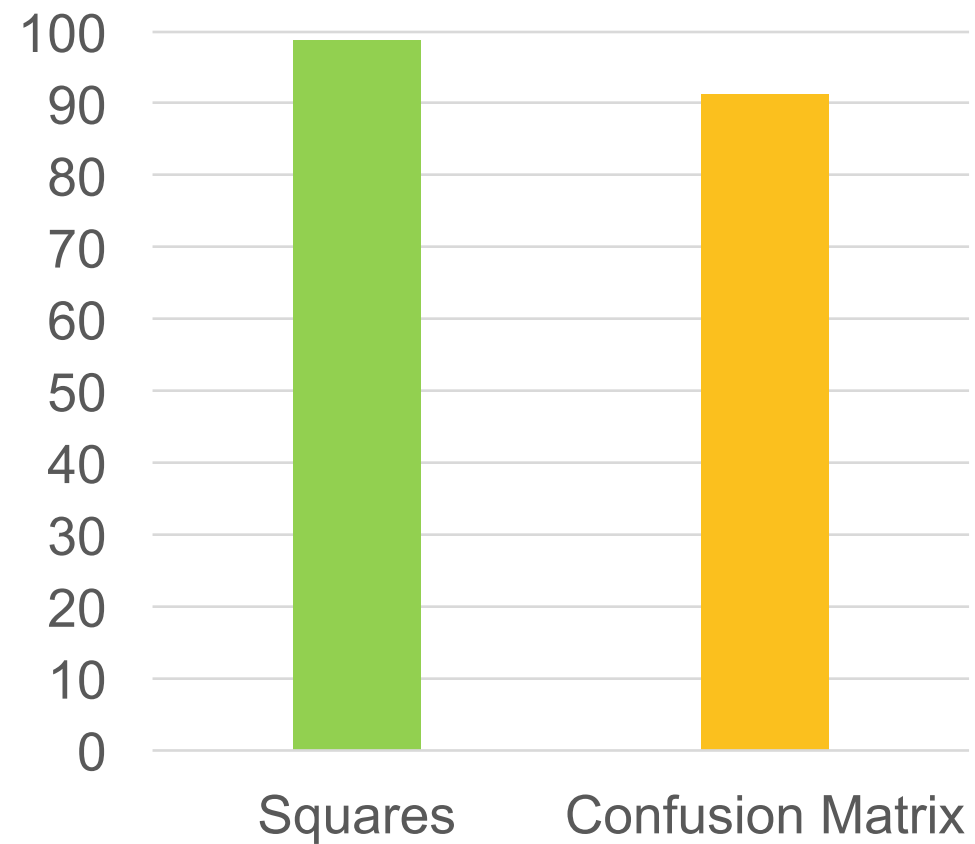Squares lead to faster task time
(Main Effect: p < 0.001)

Squares scale better in terms of the
number of classes
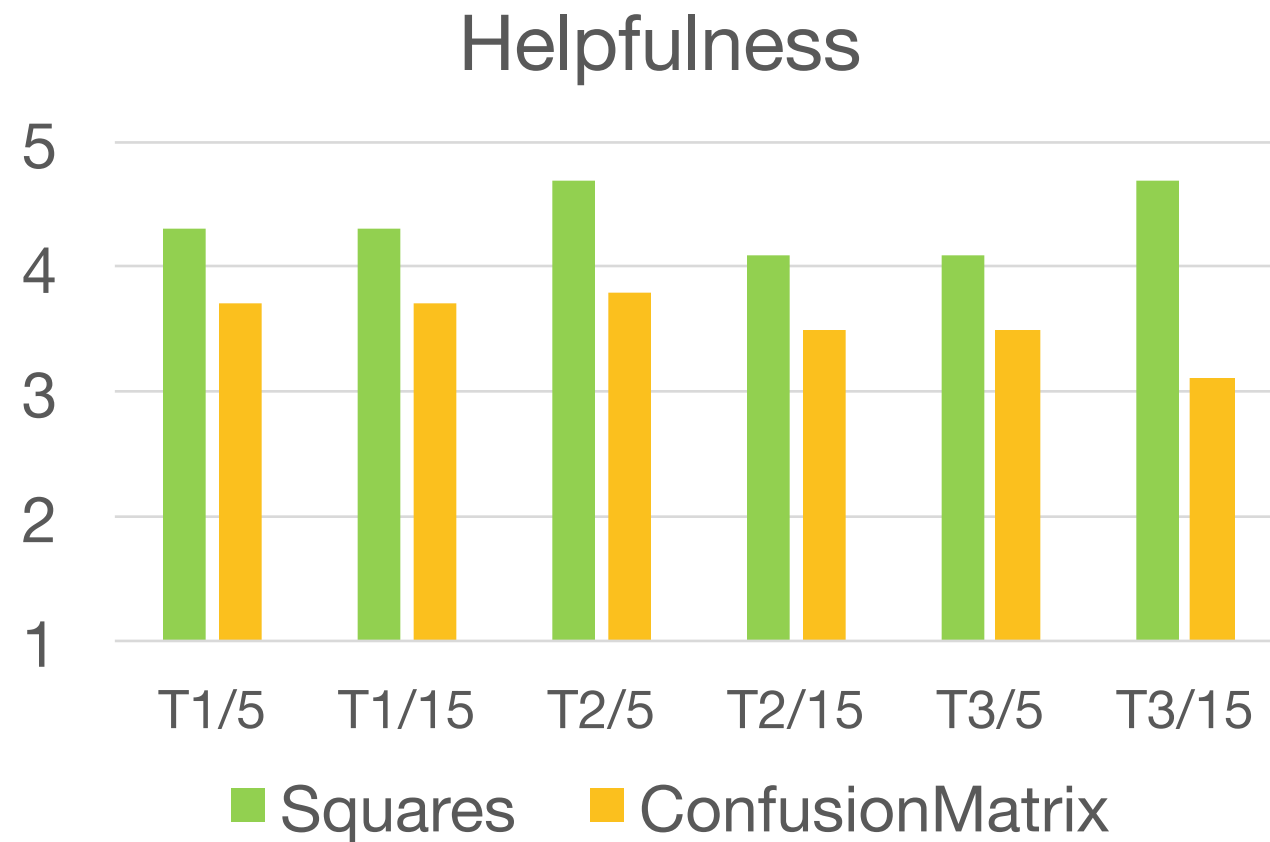(Interaction Effect: p = 0.012)
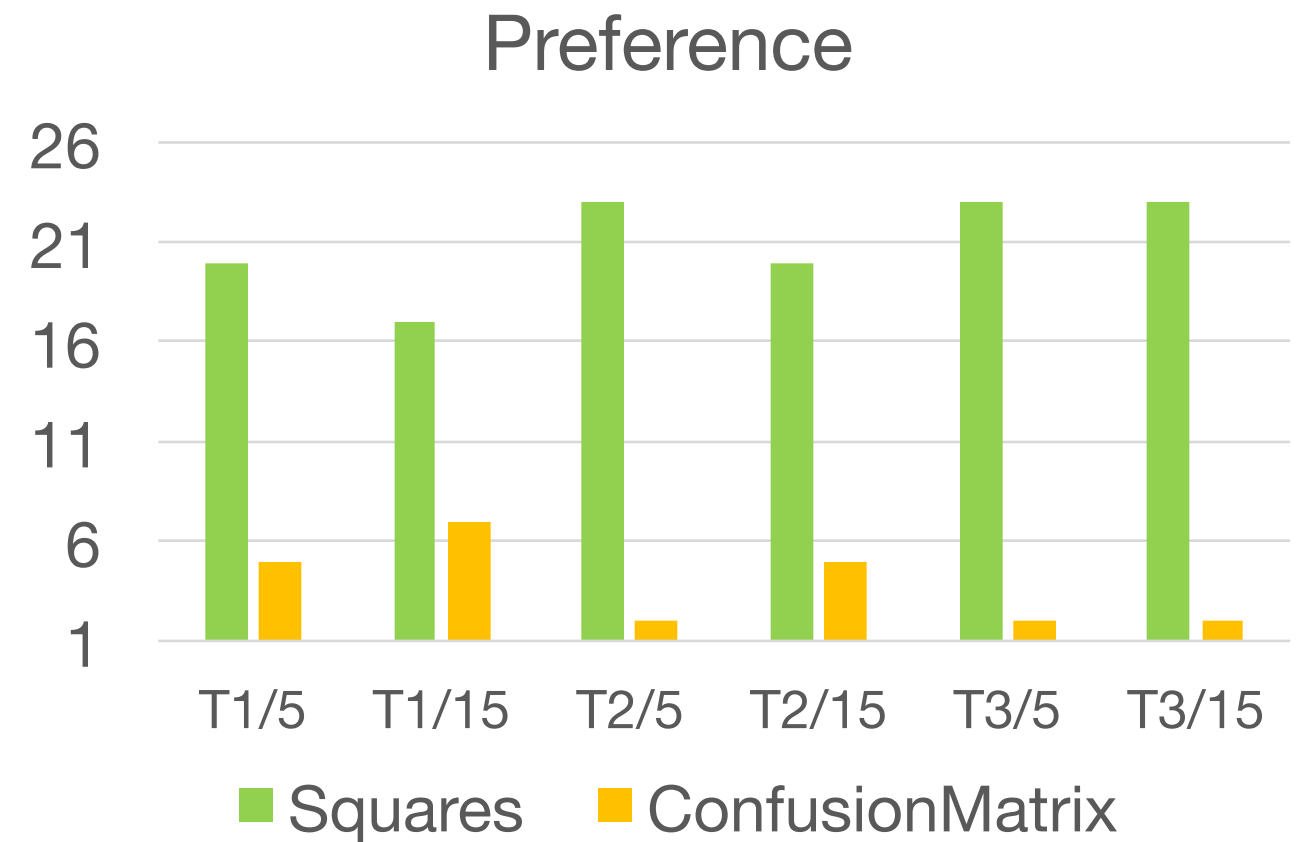
# Part 1: Squares Performed Better

- Accuracy



(p < 0.001)

- Squares lead to more accurate results

# Part 1: People Preferred Squares



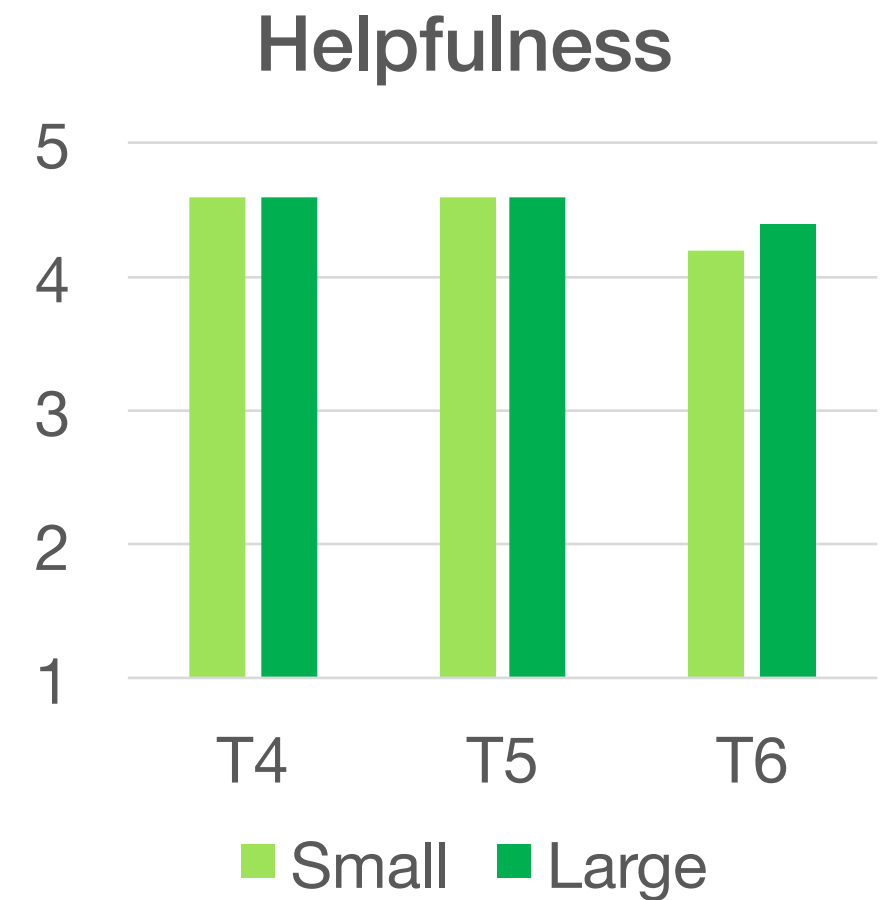Helpfulness
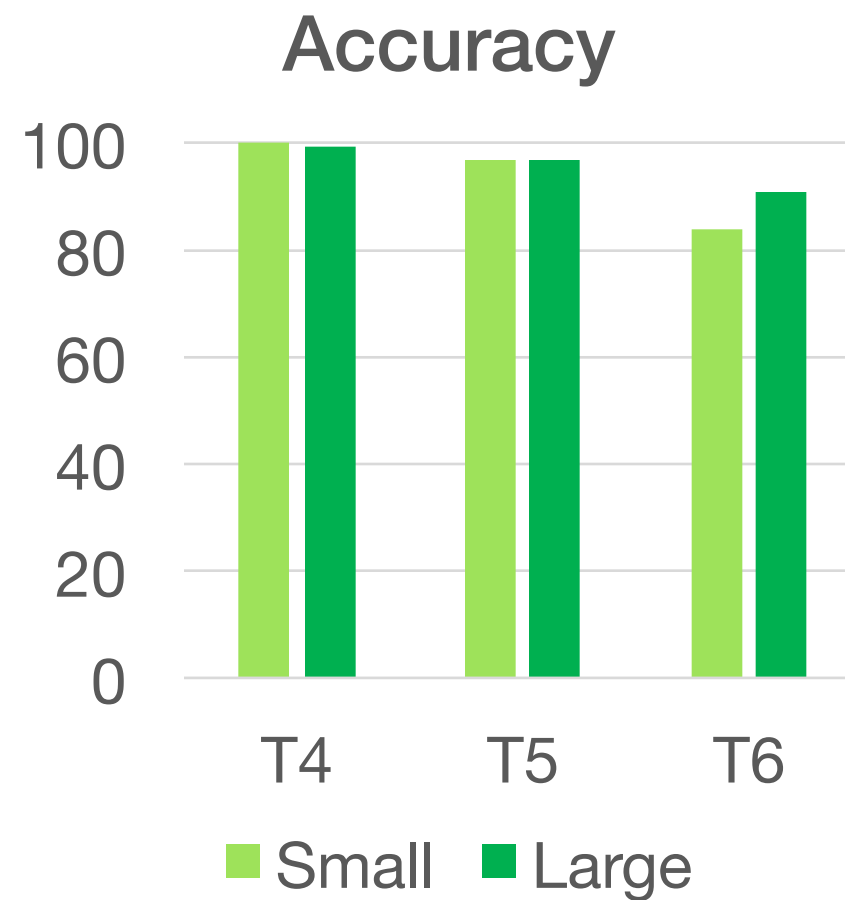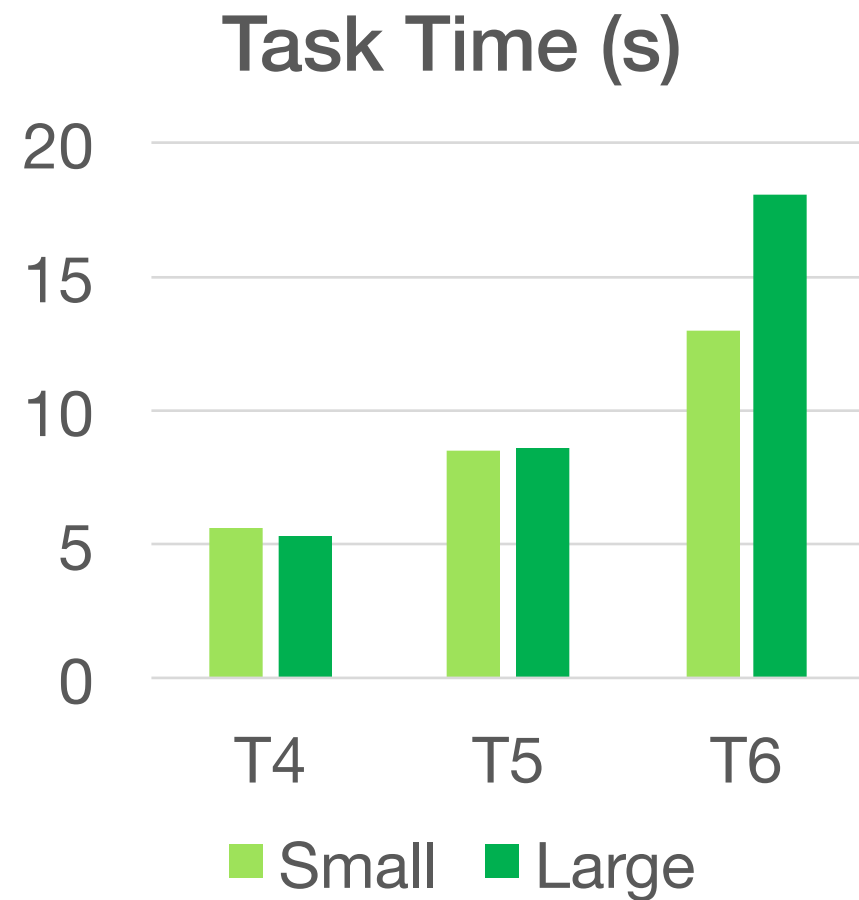
Squares was more helpful

Preference

Squares was preferred

# Part 2: (Squares Only) Distribution Tasks

- T4 – Overall
  - Select the classifier with the worst distribution

- T5 – Class-level
  - Select one of the two classes with the worst distribution

- T6 – Confusion
  - Select the two classes most confused with each other

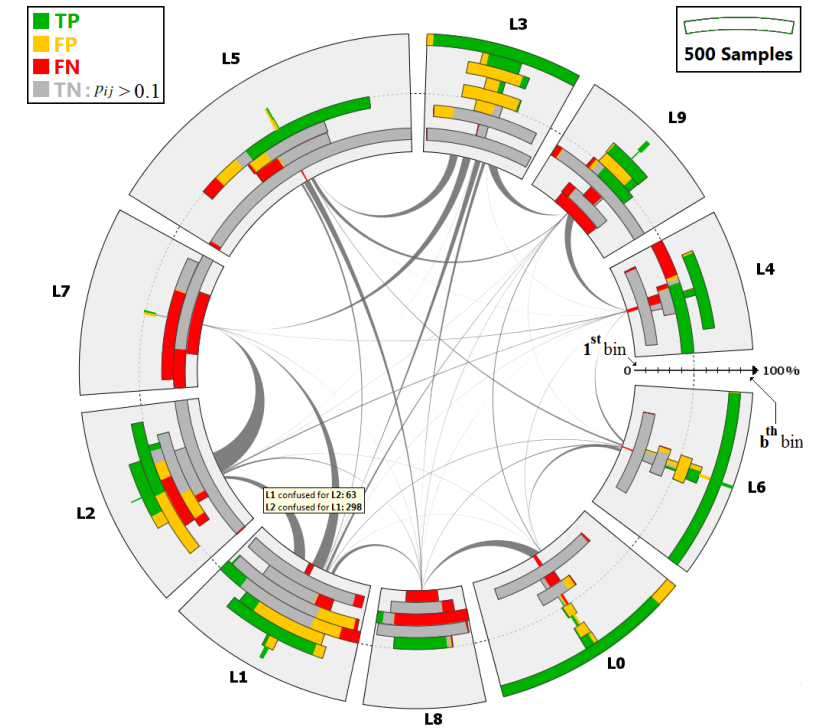# Part 2: Squares was helpful in distribution tasks

# Freeform Feedback

- Positive:
  - *"Granular and at the same time general overview of the classifiers is great."*
  - *"Seeing the distribution of scores is very helpful."*
  - *"Had fun for the first time while classifying!"*

- Negative:
  - *"I prefer having numbers than pure display."*
  - *"[Confusion Matrix is] more straightforward, lower learning curve."*

# Future Work

- ## Further Evaluation

  - Compare to alternative designs of Confusion Matrix, as well as other visualization designs in the literature

- ## Scalability

  - Supporting more than 20 classes
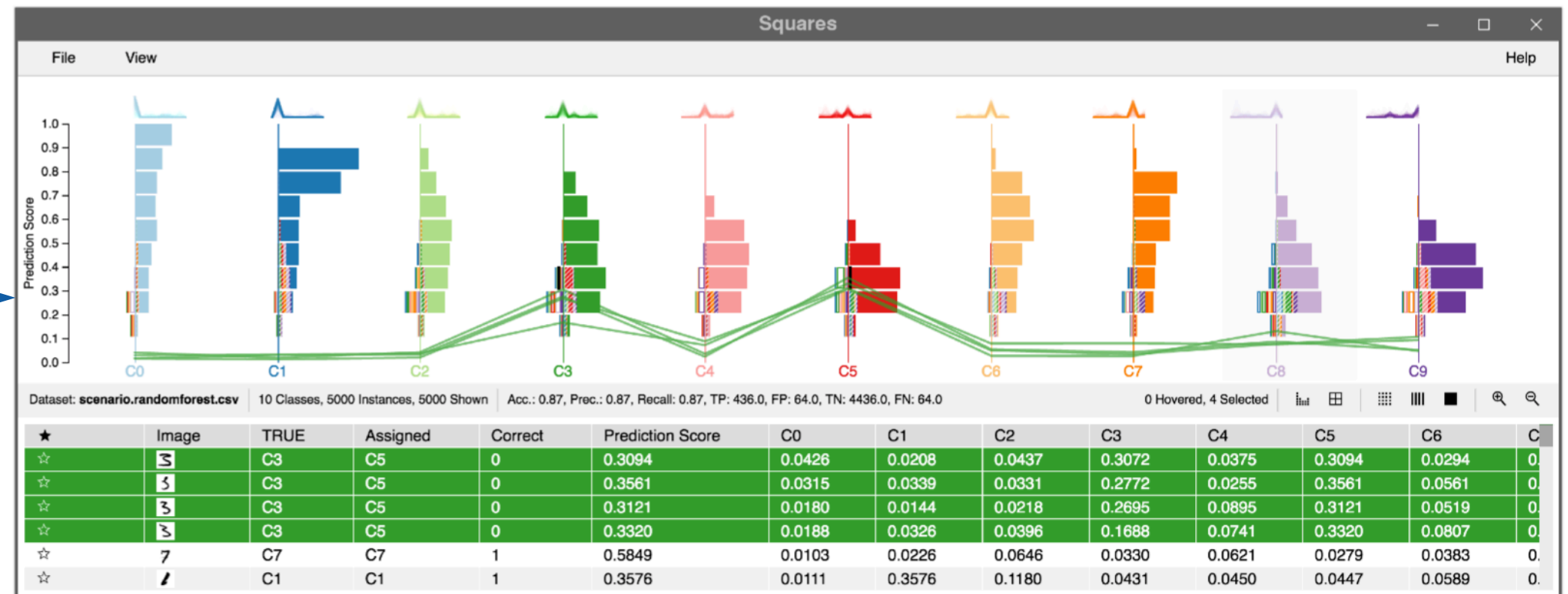  - Optimizing color assignments



Confusion Wheel [B. Alsallakh, VAST '14]
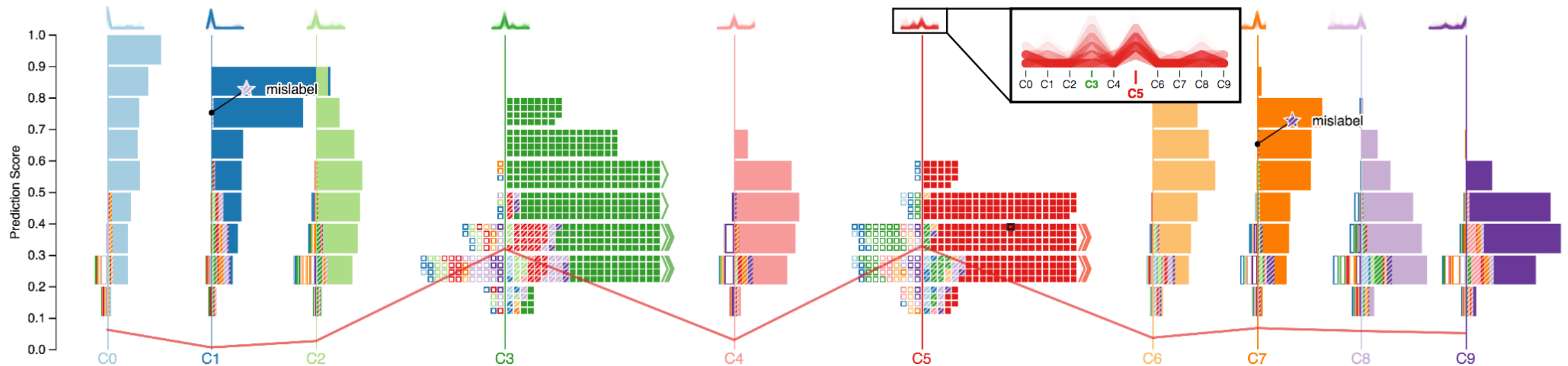
# Squares as a Tool

- Deployed along with a machine learning toolkit within Microsoft



Model Building Interface

# Acknowledgements

- We thank the support and feedback from the Machine Teaching Group at Microsoft Research.

- We thank the anonymous reviewers for their constructive comments.
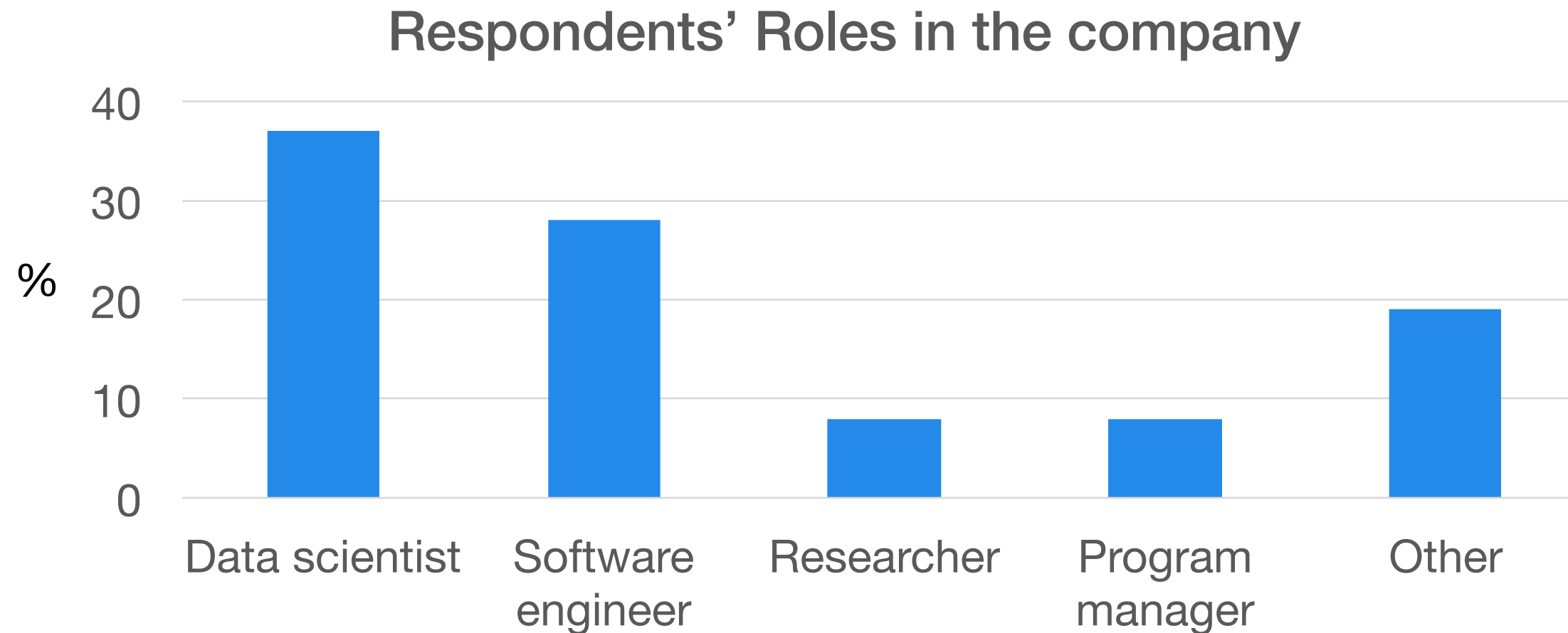
# Thanks! Questions?

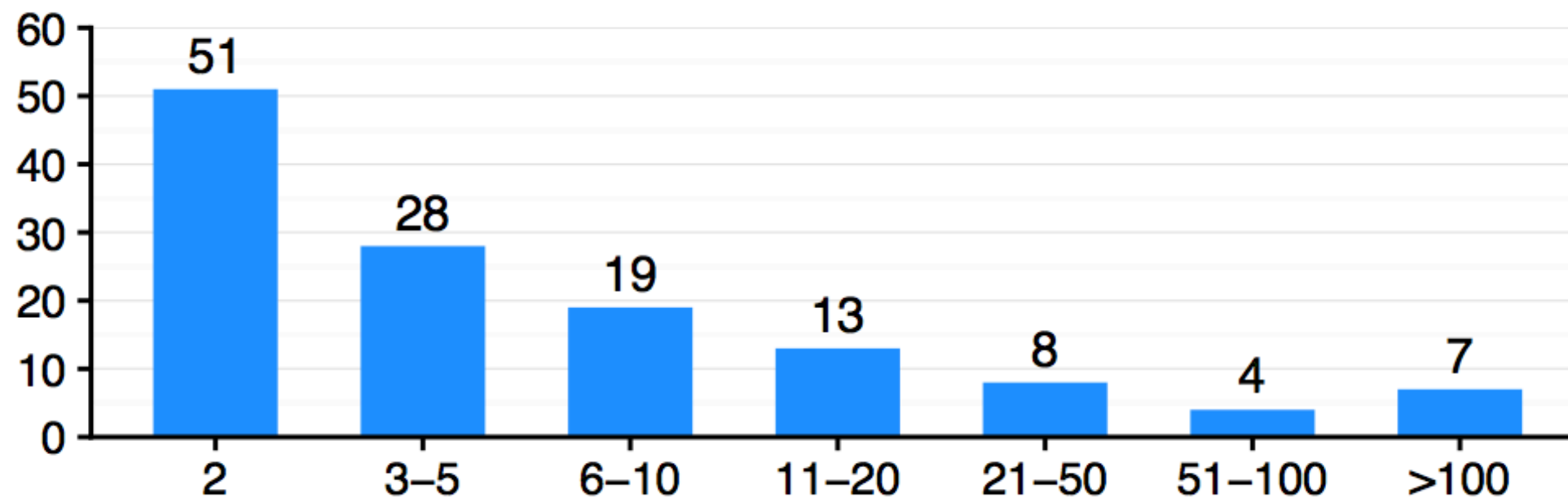Donghao Ren (donghao.ren@gmail.com)

University of California, Santa Barbara

# Survey of Machine Learning Practices

- Survey within a large software company in July. 2015.
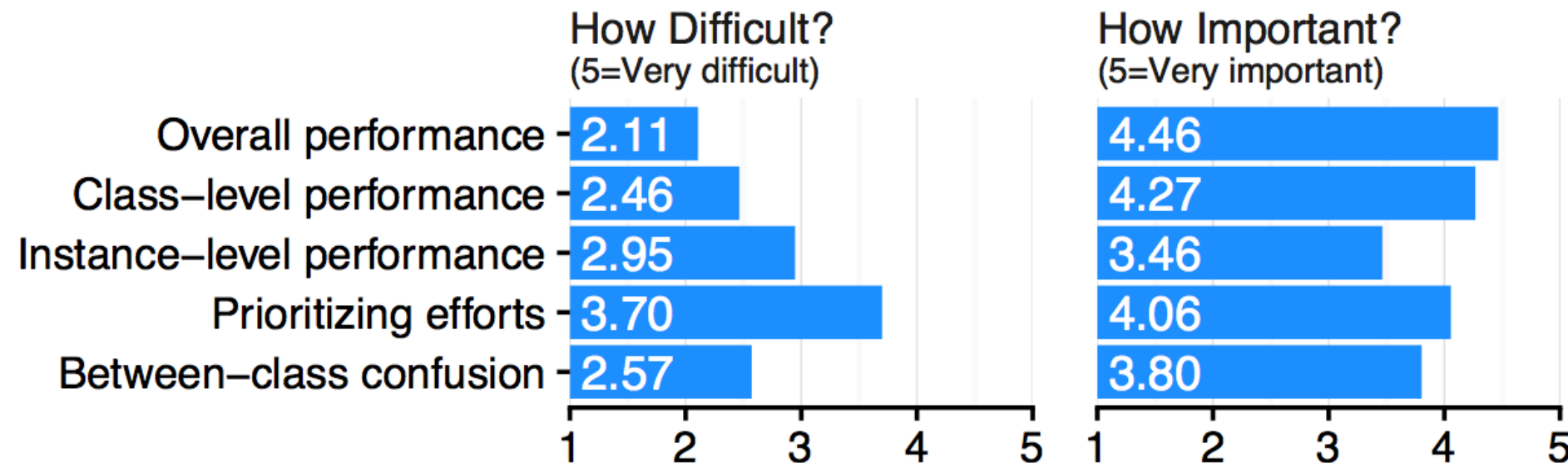- 102 respondents:

## Respondents' Roles in the company

# Number of Classes

- How many classes do your classifiers typically deal with (check all that apply)?
  - Most respondents typically deal with less than 20 classes.

# Important Tasks

- "How difficult" and "how important" ratings of tasks:
  - Prioritizing efforts is difficult even for expert users.
  - Understanding instance-level performance is relatively more difficult in common tools.

| | How Difficult? (5=Very difficult) | How Important? (5=Very important) |
|---|---|---|
| Overall performance | 2.11 | 4.46 |
| Class-level performance | 2.46 | 4.27 |
| Instance-level performance | 2.95 | 3.46 |
| Prioritizing efforts | 3.70 | 4.06 |
| Between-class confusion | 2.57 | 3.80 |

# Integrating into LUIS (Language Understanding Intelligent Service)